



---

VER MADENC L NE G R

---

# Veri Madencili ğ i Giri

---

- ğ inde ya adı ımız bili im ğ a ında elektronik ortamda mevcut verinin hızlı artı ı ve bilginin fazlala ması sebebiyle öncelikle, genelde Veri Tabanlarında Bilgi Ke fi olarak adlandırılan yeni bir paradigma ortaya ğ ıkmı tır. Daha yaygın bir kullanımla bu alana **Veri Madencili ğ i** denilmektedir.
-

# Veri Madenciliği Tanımları

(1/2)

- Veri Madenciliği (Data Mining): Büyük miktarda veri içinden, gelecekle ilgili tahmin yapmamızı sağlayacak **bağlantı** ve **kuralları** aranmasıdır. (Knowledge Discovery in Databases)
- Daha önceden bilinmeyen, geçerli ve uygulanabilir bilgilerin geniş veritabanlarından elde edilmesi ve bu bilgilerin işletme kararları verilirken kullanılmasıdır.
- Büyük ölçekli veriler arasından değerli olan bir bilgiyi elde etme işlemidir.
- Yapısal veritabanlarında depolanmış verilerden geçerli, yeni, potansiyel olarak yararlı ve nihayetinde anlaşılabilir örüntülerin tanımlanması işlemidir.

# Veri Madenciliğinin Tarihçesi (1/4)

---

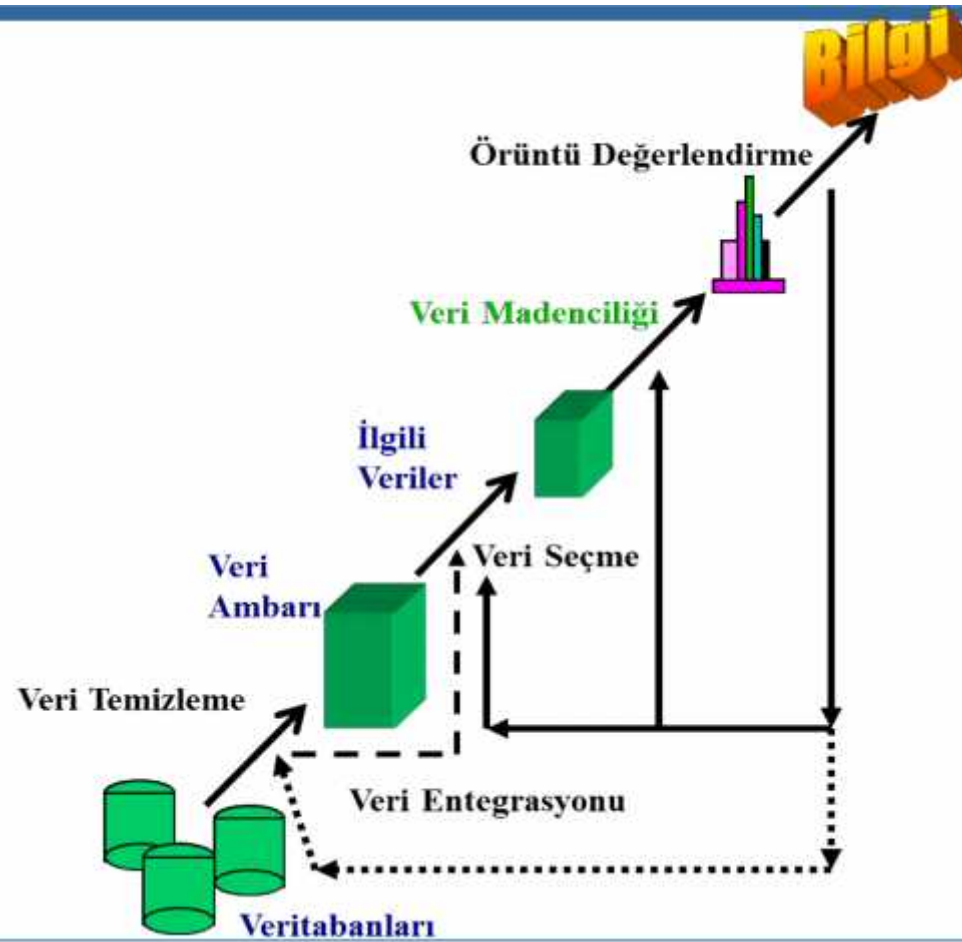
- Data FishingData Dredging: 1960
    - istatistikçiler
  - Data Mining: 1990
    - veritabanı kullanıcıları, ticari
  - Knowledge Discovery in Databases (KDD): 1989
    - Yapay zeka, makine öğrenmesi toplulukları
  - Data Archaeology, Information Harvesting, Information Discovery, Knowledge Extraction,...
-

# Bilgi Keşfi

---

- Teoride veri madenciliği bilgi keşfi ileminin amaçlarından biridir.
  - Pratikte veri madenciliği ve bilgi keşfi eş anlamlı olarak kullanılır.
  - Veri madenciliği teknikleri veriyi belli bir modele uydurur.
    - veri içindeki örüntüleri bulur
    - örüntü: veri içindeki herhangi bir yapı
  - Sorgulama ya da basit istatistik yöntemler veri madenciliği değildir.
  - Büyük veri kaynaklarından yararlı ve ilginç bilgiyi bulmak
  - Bulunan bilgi
    - gizli,
    - önemli,
    - önceden bilinmeyen,
    - yararlı olmalı.
-

# Bilgi Ke fi



# Bilgi Keşfinin Aamaları

---

- Veri Temizleme : Gürültülü ve tutarsız verileri çıkarmak
  - Veri Bütünleştirme: Birçok data kaynağını birleştirebilmek
  - Veri Seçme : Yapılacak olan analiz ile ilgili olan verileri belirlemek
  - Veri Dönüştürme : Verinin veri madenciliği yöntemine göre hale dönüştürümünü gerçekleştirmek
  - Veri Madenciliği : Verilerdeki örüntülerin belirlenmesi için veri madenciliği yöntemlerinin uygulanması
  - Örüntü Değerlendirme: Bazı ölçütlere göre elde edilmiş ilginç örüntüleri bulmak ve değerlendirmek
  - Bilgi Sunumu : Elde edilen bilgilerin kullanıcılara sunumunu
-

# Veri Madencili i Uygulama Alanları

Bilim	İş Hayatı	Web	Devlet
<ul style="list-style-type: none"><li>• Astronomi</li><li>• Biyoinformatik</li><li>• İlaç keşfi</li></ul>	<ul style="list-style-type: none"><li>• Reklam</li><li>• CRM (Müşteri İlişkileri Yönetimi) ve Müşteri Modelleme</li><li>• E-ticaret</li><li>• Yatırım değerlendirme ve karşılaştırma</li><li>• Sağlık</li><li>• Üretim</li><li>• Spor/eğlence</li><li>• Telekom (telefon ve iletişim)</li><li>• Hedef pazarlama</li></ul>	<ul style="list-style-type: none"><li>• Metin Madenciliği (haber grubu, email, dokümanlar)</li><li>• Web analizi</li><li>• Arama motorları</li></ul>	<ul style="list-style-type: none"><li>• Terörle Mücadele</li><li>• Kanun Yaptırımı</li><li>• Vergi Kaçakçılarının Profiline Çıkarılması</li></ul>



# Uygulamalar

---

- Hangi promosyonu ne zaman uygulamalıyım?
  - Hangi mü teri aldı ı krediyi geri ödemeyebilir?
  - Bir mü teriye ne kadar kredi verilebilir?
  - Sahtekarlık olabilecek davranı lar hangileridir?
  - Hangi mü teriler yakın zamanda kaybedilebilir?
  - Hangi mü terilere promosyon yapmalıyım?
  - Hangi yatırım araçlarına yatırım yapmalıyım?
-

# Veri Kaynakları

---

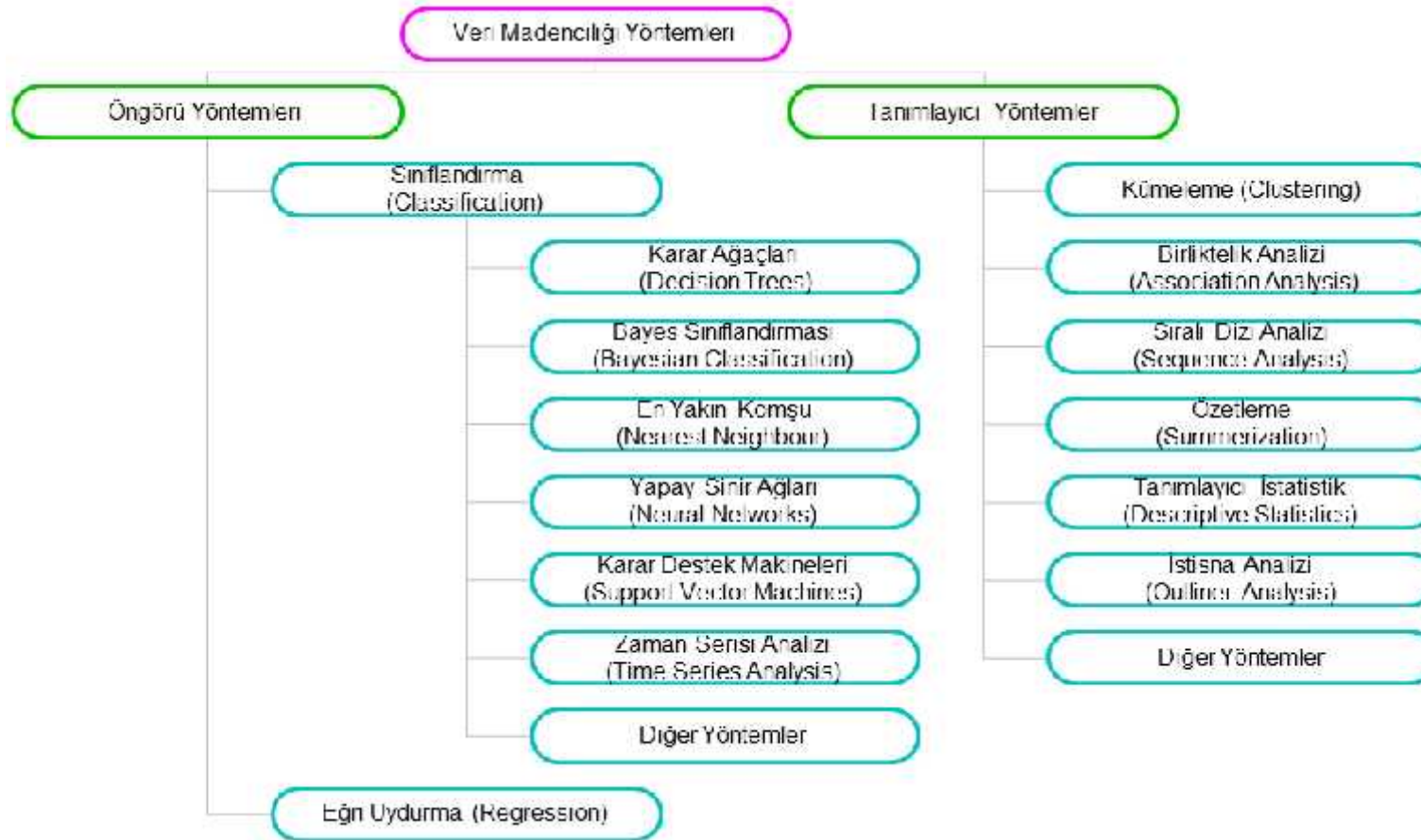
- Veri dosyaları
  - Veritabanı kaynaklı veri kümeleri
    - İlişkisel veritabanları, veri ambarları
  - Gelişen veri kümeleri
    - duraksız veri (data stream), algılayıcı verileri (sensor data)
    - zaman serileri, sıralı diziler (biyolojik veriler)
    - çizgeler, sosyal ağ (social networks) verileri
    - konumsal veriler (spatial data)
    - çoklu ortam veritabanları (multimedia databases)
    - nesneye dayalı veritabanları
    - WWW
-

# Veri Madenciliği Yöntemleri

---

- amaç: veriyi belli bir modele uydurmak
    - tanımlayıcı
      - En iyi müşteri terilerim kimler?
      - Hangi ürünler birlikte satılıyor?
      - Hangi müşteri teri gruplarının alışveriş alışkanlıkları benzer?
    - kestirime dayalı
      - Kredi başvurularını risk gruplarına ayırma
      - İşletle çalışmayı bırakacak müşteri terileri öngörme
      - Borsa tahmini
-

# Veri Madenciliği Yöntemleri



# Veri Madenciliği Seviyeleri (1/2)

---

- Sınıflandırma (Classification): Veriyi önceden belirlenmiş sınıflardan birine dahil eder.
    - Danı manlı (Gözetimli) ö renme
    - Örüntü tanıma
  - E ri uydurma (Regression): Veriyi gerçel de erli bir fonksiyona dönü türür.
  - Zaman serileri inceleme (Time Series Analysis): Zaman içinde de i en verinin de erini ö ngörür.
  - stisna Analizi (Outlier Analysis): Verinin geneline uymayan nesneleri belirleme
-

# Veri Madenciliği İlevleri (2/2)

---

- Kümeleme (Clustering): Benzer verileri aynı grupta toplama
    - Danı mansız (Gözetimsiz) öğrenme
  - Özetleme (Summarization): Veriyi alt gruplara ayırır. Her alt grubu temsil edecek özellikler bulur.
  - İli kilendirme kuralları (Association Rules)
    - Veriler arasındaki ili ki yi belirler
  - Sıralı dizileri bulma (Sequence Discovery): Veri içinde sıralı örüntüler bulmak için kullanılır.
-

# Veri Madenciliğinde Sorunlar (1/3)

- Gizlilik ve sosyal haklar
  - Kişilere ait verilerin toplanarak, kişilerden habersiz ve izinsiz olarak kullanılması
  - Veri madenciliği yöntemleri ile bulunan sonuçların izinsiz olarak açıklanması (/paylaşılması)
  - Gizlilik ve veri madenciliği politikalarının düzenlenmesi
- Kullanıcı Arabirimi
  - Görüntüleme
    - Sonucun anlaşılabilir ve yorumlanabilir hale getirilmesi
    - Bilginin sunulması
  - Etkileşim
    - Veri madenciliği ile elde edilen bilginin kullanılması
    - Veri madenciliği yöntemine müdahale etmek
    - Veri madenciliği yönteminin sonucuna müdahale etmek
- Veri madenciliği yöntemi
- Başarım ve ölçeklenebilirlik

# Veri Madenciliğinde Sorunlar (2/3)

---

- Veri madenciliği yöntemi
    - Farklı tipte veriler üzerinde çalışabilme
    - Farklı seviyelerde kullanıcı ile etkileşim halinde olabilme
    - Uygulama ortamı bilgisini kullanabilme
    - Veri madenciliği ile elde edilen sonucu anlaşılabilir şekilde sunabilme
    - Gürültülü ve eksik veri ile çalışabilme (ve iyi sonuç verebilme)
    - Değişen veya eklenen verileri kolayca kullanabilme
    - Örüntü değerlendirme: önemli örüntüleri bulma
-



# Veri Madenciliğinde Sorunlar (3/3)

---

- Başarım ve ölçeklenebilirlik
    - Kullanabilirlik ve ölçeklenebilirlik
      - Zaman karmaıklılığı ve yer karmaıklılığı kabul edilebilir
      - Örnekleme yapabilme
    - Paralel ve dağıtık yöntemler
      - Artımlı veri madenciliği
      - Parçala ve çöz
-

# VER **ÖNİŞLEME**

(Veri Öni leme-1)

---

# Veri Ön İleme

---

- Veri
  - Veri Ön İleme
    - Veriyi Tanıma
    - Veri temizleme
    - Veri birle tirme
    - Veri dönü ümü
    - Veri azaltma
  - Benzerlik ve farklılık
-

# Veri Nedir?

- Nesneler ve nesnelerin niteliklerinden oluşan küme
  - kayıt (record), varlık (entity), örnek (sample, instance) nesne için kullanılabilir.
- Nitelik (attribute) bir nesnenin (object) bir özelliğidir
  - bir insanın yaşı, ortamın sıcaklığı...
  - boyut (dimension), özellik (feature, characteristic) olarak da kullanılır.
- Nitelikler ve bu niteliklere ait değerler bir nesneyi oluşturur.

Tid	Geri Ödeme	Medeni Durum	Gelir	Dolan dancı
1	Evet	Bekar	125K	-1
2	Hayır	Evli	100K	-1
3	Hayır	Bekar	70K	-1
4	Evet	Evli	120K	-1
5	Hayır	Boşanmış	95K	1
6	Hayır	Evli	80K	-1
7	Evet	Boşanmış	220K	-1
8	Hayır	Bekar	85K	1
9	Hayır	Evli	75K	-1
10	Hayır	Bekar	90K	1

# Değer kümeleri

---

- Nitelik için saptanmış sayılar veya semboller
  - Nitelik & Değer kümeleri
    - aynı nitelik farklı değer kümelerinden değer alabilir
      - ağırlık: kg, lb(libre, ağırlık ölçüsü)
    - farklı nitelikler aynı değer kümesinden değer alabilirler
      - ID, ya : her ikisi de sayısal
-

# statistiksel Veri Türleri

- **1- Nümerik Veriler** : Sayısal-Nümerik-Nicel Veriler de denmektedir. Boy,Ya gibi süreklilik arzeden değerler Nümerik verilerdir. "Daha fazla" ifadesi ile kullanılabilirler. Sürekli ve süreksiz olarak iki başlıkta ele alınabilir:
  - a) Sürekli Nümerik Veriler: Yaş, Sıcaklık
  - b) Aralıklı Nümerik Veriler (Interval): Çocuk Sayısı, Kaza Sayısı
- **2-Nominal Veriler** : Kategorik bir veri çeşididir."Daha fazla" ifadesi ile kullanılmazlar. Kiye ayrılır:
  - a)Binary Veriler: Var-Yok, Kadın-Erkek, Hasta-Sağlıklı
  - b) kiden Çok Kategorili: Medeni Durum-Renk-Irk-Şehir, cinsiyet, Forma Numarası
  - Örneğin forma numarası oyuncunun seviyesi ile ilgili bir bilgi içermez.

# statistiksel Veri Türleri

---

- **3-Ordinal Veriler** : Ordinal veriler de yine kategorik veri türündendir. Fakat de erleri arasında sıralı bir ili ki bulunmaktadır. “Daha fazla” ifadesi ile kullanılabilirler ancak ne kadar daha fazla oldu unun ölçüsünü veremezler. Örne im: E itim Düzeyi, Sosyoekonomik ölçek skorları gibi. Nominal veriler, ordinal verilere göre daha az bilgi ta ırlar.
  - **4-Ratio Veriler** : Nümerik verilere benzerler. 100 santigrat derece, 50 santrigat derecenin iki katı denilemez ama derece kelveine çevrilirse 60 kelvin 30 kelvinin 2 misli sıcak denilebilir. Oran verilebilir veri türlerine Ratio veriler denir. Burada kelvin derece ratio türünden bir de i ken iken, santigrat ise nümerik veri türüne örnek olarak verilebilir.
-

# Nitelik Türleri

---

- Belli aralıkta yeralan de i kenler (interval)
    - sıcaklık, tarih
  - kili de i kenler (binary)
    - cinsiyet
  - Ayırık ve sıralı de i kenler
    - göz rengi, posta kodu
-



# Problem

- Gerçek uygulamalarda toplanan veri kirli
  - eksik: bazı nitelik de erleri bazı nesneler için girilmemi , veri madencili i uygulaması için gerekli bir nitelik kaydedilmemi
    - meslek = " "
  - gürültülü: hatalar var
    - maa = "-10"
  - tutarsız: nitelik de erleri veya nitelik isimleri uyumsuz
    - ya = "35", d.tarihi: "03/10/2004"
    - önceki oylama de erleri: "1,2,3", yeni oylama de erleri: "A,B,C"
    - bir kaynakta nitelik de eri 'ad', di erinde 'isim'

# Veri kirliliği örneği-1

kapsam	sorun	Kirli veriler	sebep
özellik	Yanlış değer	Doğum_günü =30.13.1990	Değerler alan dışındadır
Kayıt	Özellikler arasında bağımlılığın yanlış olması	Yaş=42 Doğum_günü=12.02.1990	«yaş»la doğum günü değerleri tutarsızdır
Kayıt türü	Eşsizliğin bozulması	Pers1=(ad=«Ali Yavuz», pno=«123456»  Pers1=(ad=«Metin SAĞLAM», pno=«123456»	Personel numarasının eşsiz olması koşulu bozulmuştur
kaynak	Erişimsel bütünlüğün bozulması	Pers1=(ad=«Metin SAĞLAM», şube_no=«123456»	«123456»no'lu şube tanımlanmamıştır

# Veri kirliliği örneği

kapsam	sorun	Kirli veriler	sebebi
özellik	Değer yoktur	<u>Tel: =285218</u> 163	Rakam eksiktir
özellik	Kelimenin yanlış yazılışı	Kent=«Trabzun»	Fonetik hata
özellik	yanlış alan değeri	Kent=«İtalya»	«İtalya» «kent» alanına dahil değil
kayıt	Özellikler arası bağımlılığın bozulması	Kent=«Çanakkale»; plaka_no=19	«Çanakkale'nin plaka numarası 19 değil
Kayıt türü	Kelimelerin farklı dizilişi	Ad1 =«Kerim UĞUR» Ad2=«YILMAZ Temel»	Ad ve soyadların sıraları farklıdır
Kayıt türü	Kayıtlarda zıtlık	Pers1=(ad=«Ali Yavuz», doğum_tar=12.12.1995  Pers2=(ad=«Ali Yavuz», doğum_tar=10.09.1995	Aynı varlık farklı değerlerle tanımlanmıştır

# Verinin Kirli Olma Nedenleri

---

- Eksik veri kayıtlarının nedenleri
    - Veri toplandı ı sırada bir nitelik de erinin elde edilememesi, bilinmemesi
    - Veri toplandı ı sırada bazı niteliklerin gereklili inin görülememesi
    - nsan, yazılım ya da donanım problemleri
  - Gürültülü (hatalı) veri kayıtlarının nedenleri
    - Hatalı veri toplama gereçleri
    - nsan, yazılım ya da donanım problemleri
    - Veri iletimi sırasında problemler
  - Tutarsız veri kayıtlarının nedenleri
    - Verinin farklı veri kaynaklarında tutulması
    - levsel ba ımlılık kurallarına uyulmaması
-

# Sonuç

---

- Veri güvenilirmez oldu unda:
    - Veri madencili i sonuçlarına güvenilebilir mi?
    - Kullanılabilir veri madencili i sonuçları kaliteli veri ile elde edilebilir.
  - Veri kaliteli ise veri madencili i uygulamaları ile yararlı bilgi bulma ansı daha fazla.
-



---

# Veriyi Tanıma

---

# Veriyi Tanımlayıcı Özellikler

---

- Amaç: Veriyi daha iyi anlamak
    - Merkezi eğilim (central tendency), varyasyon, yayılma, dağılım
  - Verinin dağılım özellikleri
    - Ortanca, en büyük, en küçük, sıklık derecesi, aykırılık, varyans
  - Sayısal nitelikler -> sıralanabilir değerler
    - verinin dağılımı
    - kutu grafiği çizimi ve sıklık derecesi incelemesi
-

# Veri Ön İşleme

---

- Veri temizleme
    - Eksik nitelik değerlerini tamamlama, hatalı veriyi düzeltme, aykırılıkları saptama ve temizleme, tutarsızlıkları giderme
  - Veri birleştirme
    - Farklı veri kaynaklarındaki verileri birleştirme
  - Veri dönüştürme
    - Normalizasyon
  - Veri azaltma
    - Aynı veri madenciliği sonuçları elde edilecek şekilde veri miktarını azaltma
-



# Veri Temizleme

---

- Gerçek uygulamalarda veri eksik, gürültülü veya tutarsız olabilir.
  - Veri temizleme i lemleri
    - Eksik nitelik de erlerini tamamlama
    - Aykırılıkların bulunması ve gürültülü verinin düzeltilmesi
    - Tutarsızlıkların giderilmesi
-

# Eksik Veri

---

- Veri için bazı niteliklerin de erleri her zaman bilinemeyebilir.
  - Eksik veri
    - di er veri kayıtlarıyla tutarsızlı ı nedeniyle silinmesi
    - bazı nitelik de erleri hatalı olması dolayısıyla silinmesi
    - yanlış anlama sonucu kaydedilmeme
    - veri giri i sırasında bazı nitelikleri önemsiz görme
-

# Eksik Veriler nasıl Tamamlanır?

---

- Eksik nitelik değerleri olan veri kayıtlarını kullanma (sil)
  - Eksik nitelik değerlerini elle doldur
  - Eksik nitelik değerleri için global bir değeri kullan (Null, bilinmiyor,...)
  - Eksik nitelik değerlerini o niteliğin ortalama değeri ile doldur
  - Aynı sınıfa ait kayıtların nitelik değerlerinin ortalaması ile doldur
  - Olasılığı en fazla olan nitelik değeriyle doldur
-

Age	Income	Team	Gender
23	24,200	Red Sox	M
39	?	Yankees	F
45	45,390	?	F

Olasılığı en fazla olan nitelik değeriyle doldur

Eksik nitelik değerlerini o niteliğin ortalama değeri ile doldur

Aynı sınıfa ait kayıtların nitelik değerlerinin ortalaması ile doldur

# Gürültülü Veri

---

- Ölçülen bir de erdeki hata
  - Yanlı nitelik de erleri
    - hatalı veri toplama gereçleri
    - veri giri i problemleri
    - veri iletimi problemleri
    - teknolojik kısıtlar
    - nitelik isimlerinde tutarsızlık
-

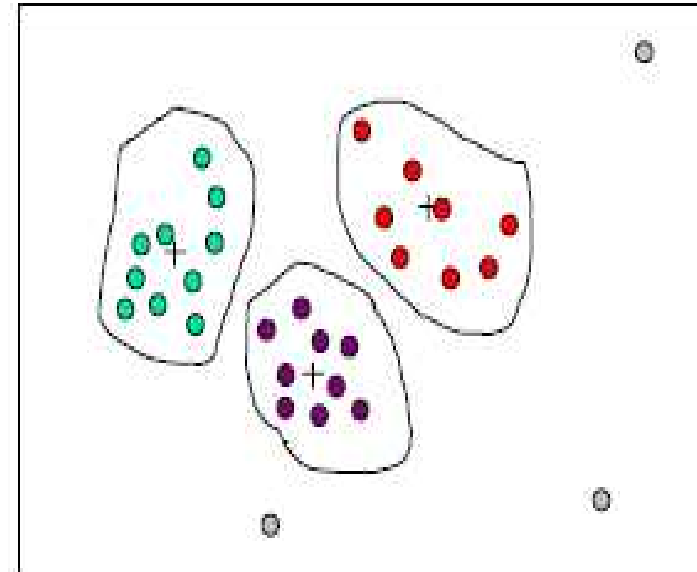
# Gürültülü Veri nasıl düzeltilir?

---

- Gürültüyü yok etme
    - Kümeleme
    - E ri uydurma
      - veriyi bir fonksiyona uydurarak gürültüyü düzeltir.
-

# Kümeleme

- Benzer veriler aynı kümede olacak şekilde gruplanır
- Bu kümelerin dışında kalan veriler aykırılık olarak belirlenir ve silinir.



## E ri Uydurma

---

- Veri bir fonksiyona uydurulur. Doğrusal veri uydurmada, bir de i kenin de eri di er bir de i ken kullanılarak bulunabilir.
-





---

Veri Birle tirme

---

# Veri Birle tirme

---

- Farklı kaynaklardan verilerin tutarlı olarak birle tirilmesi
    - ema birle tirilmesi
      - Aynı varlıkların saptanması
  - Nitelik de erlerinin tutarsızlı ının saptanması
    - Aynı nitelik için farklı kaynaklarda farklı de erler olması
    - Farklı metrikler kullanılması
-

# Gereksiz Veri

---

Farklı veri kaynaklarından veriler birleştirilince gereksiz (fazla) veri oluşabilir

- aynı nitelik farklı kaynaklarda farklı isimle
  - bir niteliğin değeri başka bir nitelik kullanılarak hesaplanabilir
-

# VERİ MADENCİLİK

(Veri Ön İşleme-2)

---

# Veri Dönü üümü

---

- Veri, veri madencili i uygulamaları için uygun olmayabilir
- Seçilen algoritmaya uygun olmayabilir
  - Veri belirleyici de il

## Çözüm

- Veri düzeltme
  - Normalizasyon
-

# Normalizasyon

- min-max normalizasyon
  - min-max normalle tirmesi ile orijinal veriler yeni veri aralığına dönüşüm ile dönüştürülürler. Bu veri aralığı genellikle 0-1 aralığıdır.
- z-score normalizasyon
  - z Skoru normalle tirmede (veya 0 ortalama normalle tirme) ise de i kenin her hangi bir y de eri, de i kenin ortalaması ve standart sapmasına bağlı olarak bilinen Z dönüşümü ile normalle tirilir.
- ondalık normalizasyon
  - Ondalık ölçekleme ile normalle tirmede ise, ele alınan de i kenin de erlerinin ondalık kısmı hareket ettirilerek normalle tirme gerçekleştirilir. Hareket edecek ondalık nokta sayısı, de i kenin maksimum mutlak de erine bağlıdır. Ondalık ölçeklemenin formülü aşağıdaki gibidir:
    - Örneğin 900 maksimum de er ise,  $n=3$  olacağından 900 sayısı 0,9 olarak normalle tirilir.

# Normalizasyon

- min-max normalizasyon

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new\_max}_A - \text{new\_min}_A) + \text{new\_min}_A$$

- z-score normalizasyon

$$v' = \frac{v - \text{mean}_A}{\text{stand\_dev}_A}$$

- ondalık normalizasyon

$$v' = \frac{v}{10^j} \quad j: \text{Max}(|v'|) < 1 \text{ olacak \u015fekildeki en k\u00fc\u00e7\u00fck tam say\u0131}$$

# Normalizasyon

- **Min-max normalizasyon:**

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new\_max}_A - \text{new\_min}_A) + \text{new\_min}_A$$

- Ör. Yıllık gelir \$12,000 ile \$98,000 arasını [0.0, 1.0] aralığına normalize edelim. \$73,000 kaç denge gelir?

$$\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$$

- **Z-score normalizasyon** ( $\mu$ : ortalama,  $\sigma$ : standard sapma):

$$v' = \frac{v - \mu_A}{\sigma_A}$$

- Ör. Let  $\mu = 54,000$ ,  $\sigma = 16,000$ . Öyleyse:  $\frac{73,600 - 54,000}{16,000} = 1.225$

- **Ondalıklı Normalizasyon**

$$v' = \frac{v}{10^j} \quad \$73,000 \text{ kaç denge gelir? } v' = 0.73$$



# Nitelik Olu turma

---

- Yeni nitelikler yarat
    - orjinal niteliklerden daha önemli bilgi içersin
      - alan=boy x en
    - veri madencili i algoritmalarının ba arımı daha iyi olsun
-



---

# Veri Azaltma

---

# Veri Azaltma

---

- Veri miktarı çok fazla oldu u zaman veri madencili i algoritmalarının çalışması ve sonuç üretmesi çok uzun sürebilir
    - veriyi azaltma başarımlarını artırır
    - sonucun (nerdeyse) hiç de i memesi gerekir
  - Veri azaltma
    - nitelik azaltma
    - veri sıkı tırma
    - veri ayırıkla tırma
    - veri küçültme
-

# Nitelik Azaltma

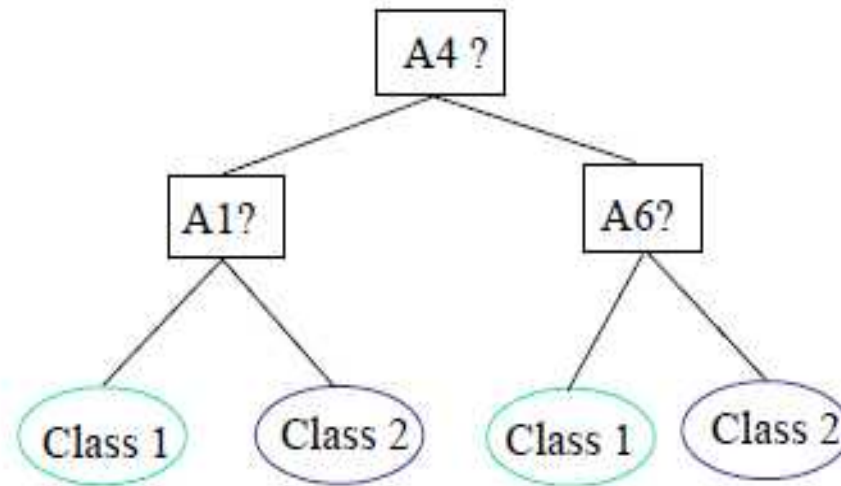
---

- Nitelikler kümesinin bir alt kümesi seçilerek veri madenciliği için kullanılır.
  - $d$  boyutlu veri kümesi  $k < d$  olacak şekilde  $k$  boyuta indirilir.
  - Nitelik seçme
    - Veri madenciliği uygulaması için gerekli olan niteliklerin seçilmesi
-

# Örnek

Başlangıç nitelikler kümesi:

$\{A1, A2, A3, A4, A5, A6\}$



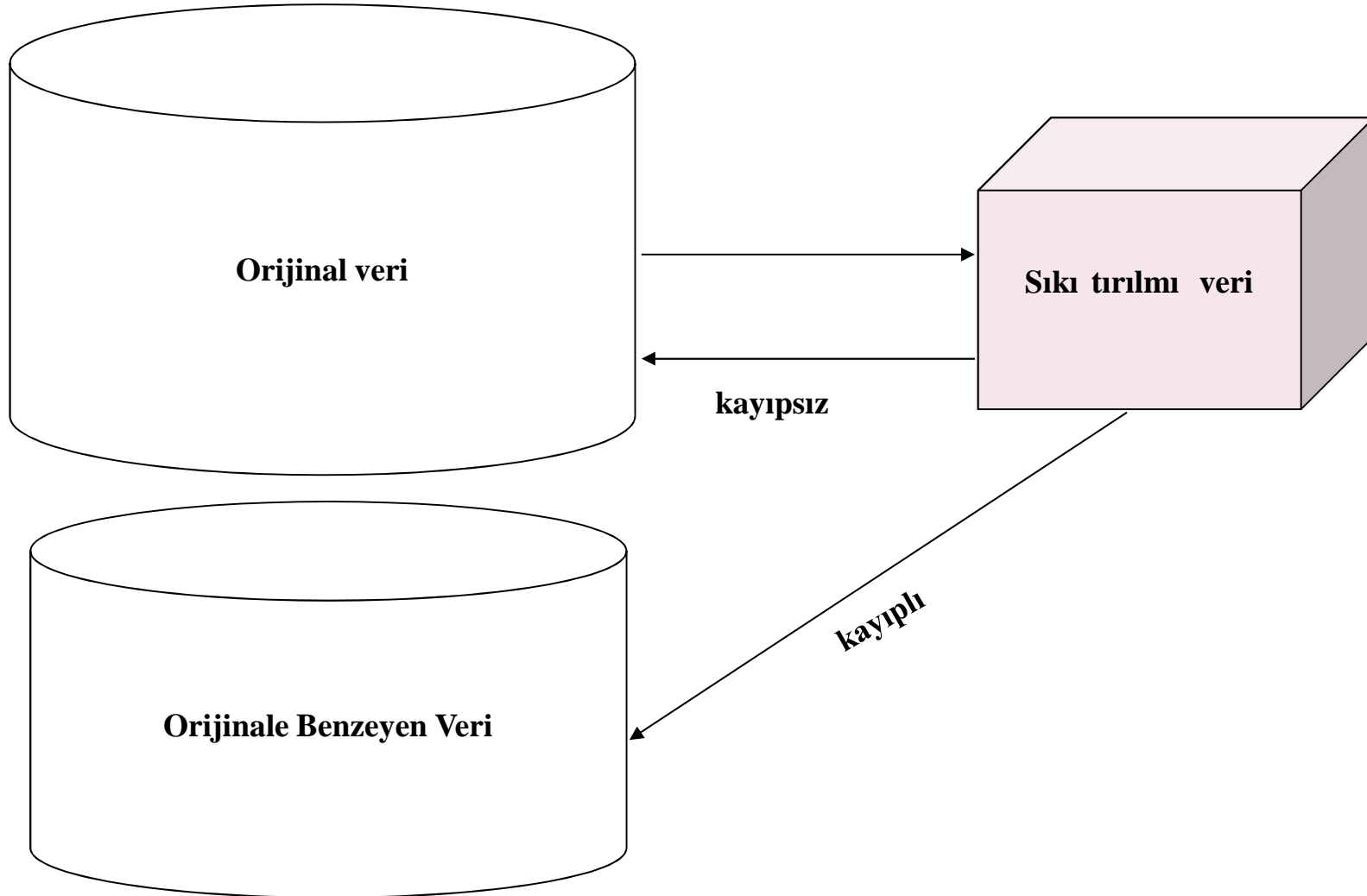
Seçilen nitelik kümesi:  $\{A1, A4, A6\}$

# Veri Sıkıştırma

---

- Verinin boyutunu azaltır
    - daha az saklama ortamı
    - veriye ulaşmak daha hızlı
  - Kayıplı ve kayıpsız veri sıkıştırma
    - bazı yöntemler bazı veri tiplerine uygun
  - Eğer veri madenciliği yöntemi sıkıştırılmış veri üzerinde doğrudan çalışabiliyorsa elverişli
-

# Veri Sıkıştırma



# Veri Ayırıkla tırma

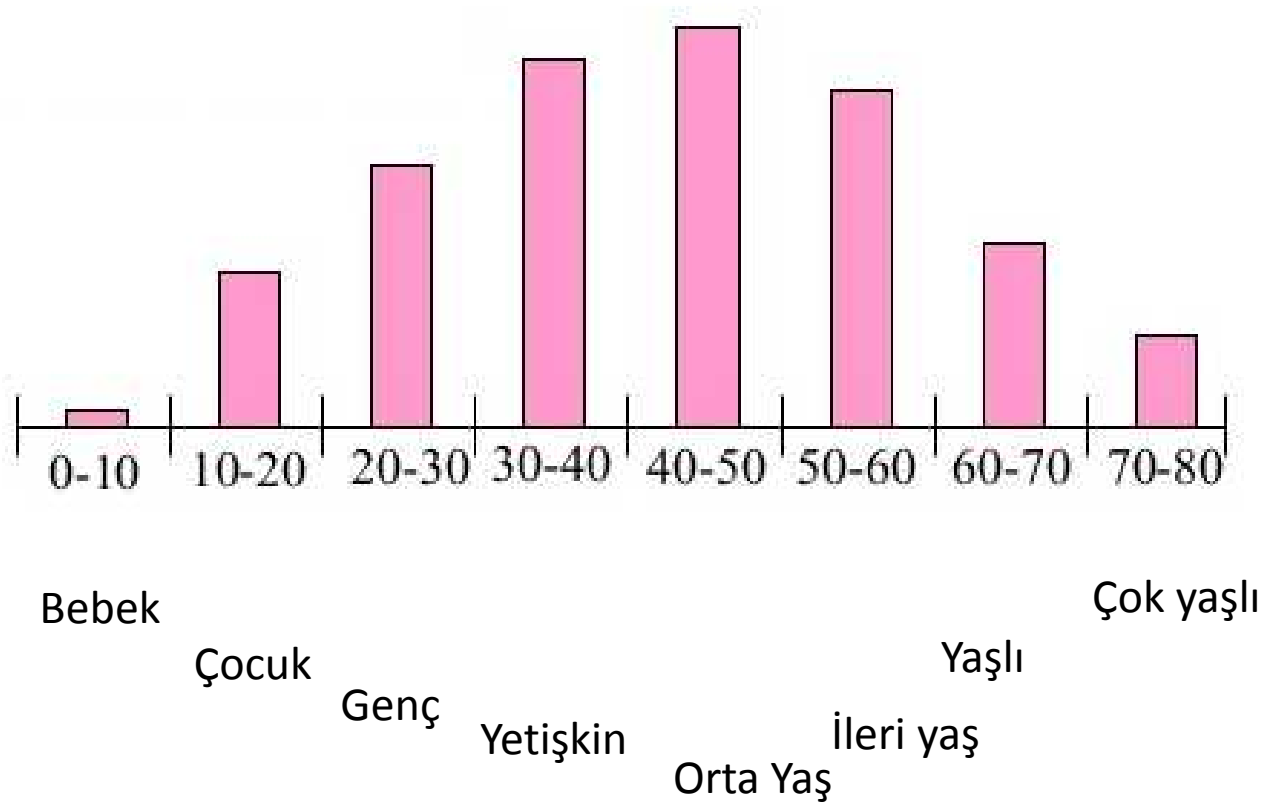
---

- Bazı veri madencili i algoritmaları sadece ayırık veriler ile alı ır.
  - Sürekli bir nitelik de erini bölerek her aralı ı etiketler.
  - Verinin de eri, bulundu u aralı ın etiketi ile de i ir.
  - Veri boyutu küçölür.
-



# Veri Ayırıkla tırma

Müşteri Yaşına göre ayırıklaştırma



# Veri Küçültme

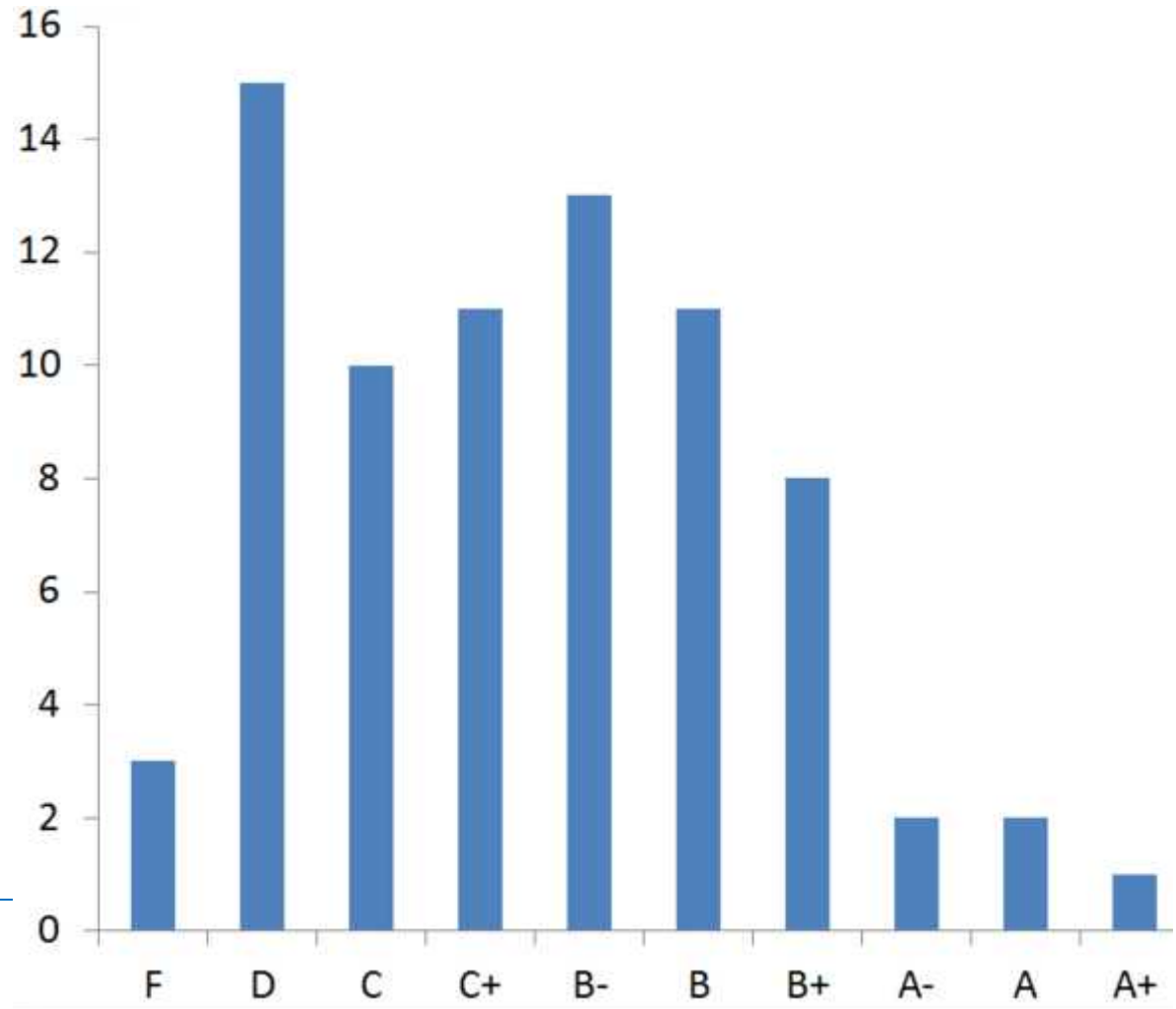
---

## ■ Veriyi farklı şekillerde gösterme

- histogram
  - kümeleme
  - örnekleme
-

# Histogram ile Veri Küçültme

- Verinin dağılımı
- Veriyi bölerek her bölüm için veri de erini gösterir (toplam, ortalama)



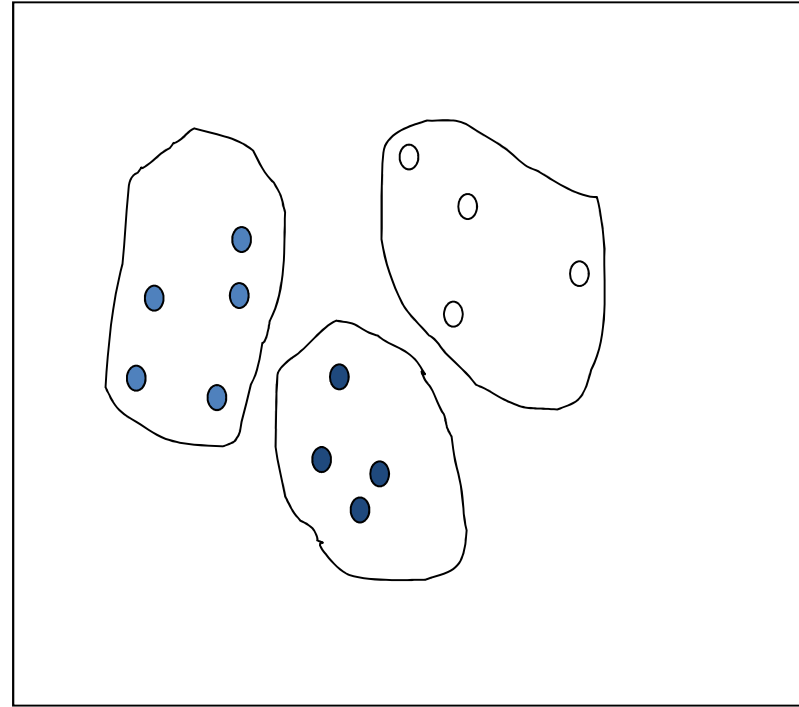
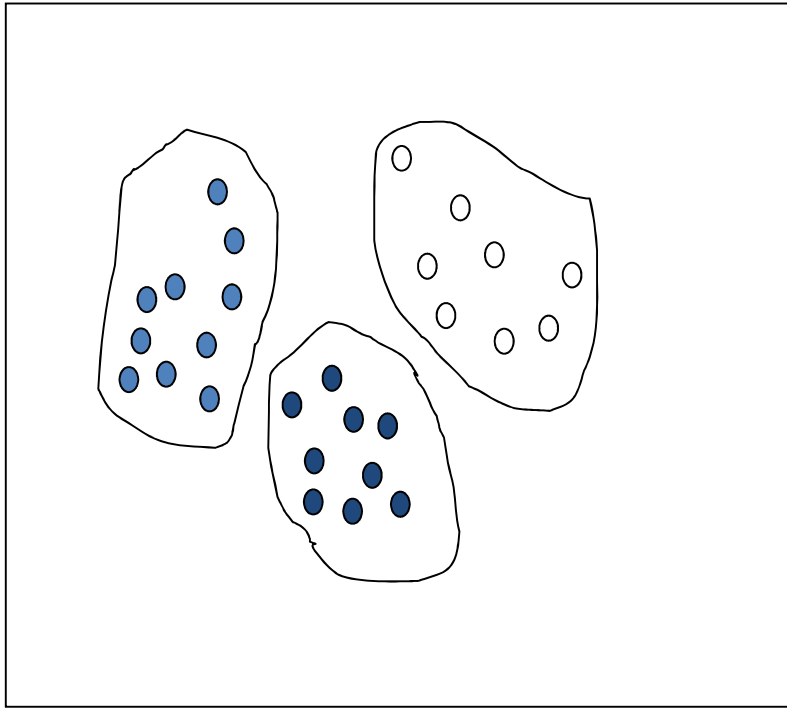
# Kümeleme ile Veri Küçültme

---

- Veri kümelerine ayrılır
  - Veri kümeleri temsil eden örnekler (küme merkezleri) ve aykırılıklar ile temsil edilir
  - Etkisi verinin dağılımına bağlı.
-

# Kümeleme ile veri küçültme

- Kümelenmiş veri
- Her kümeden orantılı sayıda temsilci seçimi

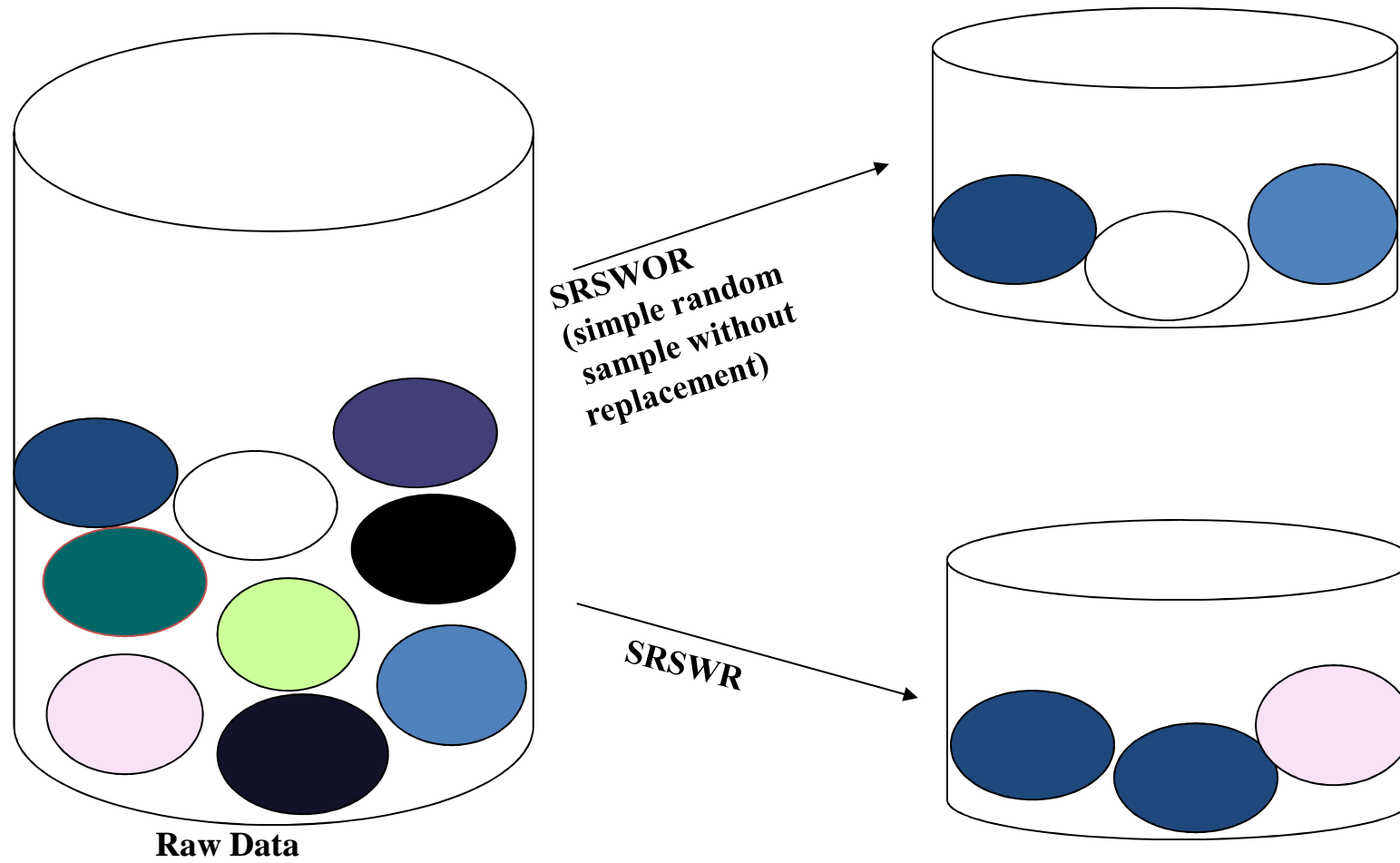


# Örnekleme ile Veri Küçültme

---

- Büyük veri kümesini daha küçük bir alt küme ile temsil etme
  - Alt küme nasıl seçiliyor?
    - yerine koymadan örnekleme (SRSWOR)
    - yerine koyarak örnekleme (SRSWR)
    - katman örnekleme (katman: nitelik de erine göre grup)
-

# Örnekleme





---

# Benzerlik ve Farklılık

---



# Benzerlik ve Farklılık

---

## ■ Benzerlik

- iki nesnenin benzerliğini ölçen sayısal değer
- nesneler birbirine daha benzer ise daha büyük
- genelde 0-1 aralığında değer alır

## ■ Farklılık

- iki nesnenin birbirinden ne kadar farklı olduğunu gösteren sayısal değer
  - nesneler birbirine daha benzer ise daha küçük
  - en küçük farklılık genelde 0
  - üst sınır değişebilir.
-

# Uzaklık Çe itleri

---

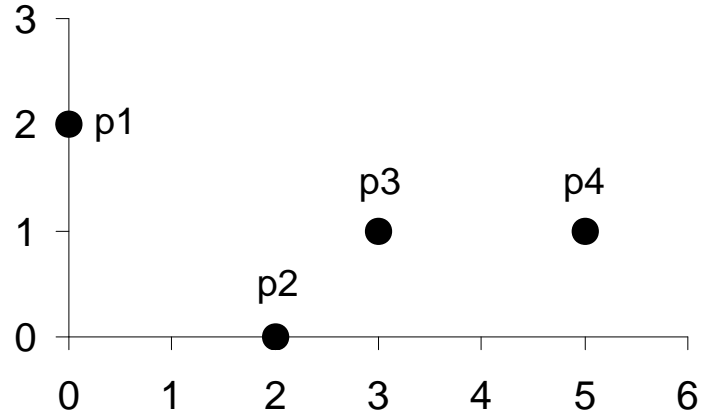
- Öklid
  - Minkowski (Manhattan)
-

# Öklid Uzaklığı

- Öklid uzaklığı (Euclidean Distance) nesneler arası farklılığı bulmak için kullanılır.
  - $p$  adet niteliği (boyutu) olan  $i$  ve  $j$  nesneleri arasındaki uzaklık

$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2}$$

# Öklid Uzaklığı



point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

**Uzaklık Matrisi**

# Minkowski Uzaklığı

- Öklid uzaklığının genelleştirilmiş hali

$$d(i, j) = \sqrt[q]{(|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + \dots + |x_{ip} - x_{jp}|^q)} \quad q: \text{pozitif tam sayı}$$

- $q=1 \rightarrow$  Manhattan uzaklığı
-

# Minkowski Uzaklığı

## Manhattan Uzaklık Matrisi

point	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

L1	p1	p2	p3	p4
p1	0	4	4	6
p2	4	0	2	4
p3	4	2	0	2
p4	6	4	2	0

## Öklid Uzaklık Matrisi

L2	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

# Benzerlik Özellikleri

---

- ki nesne arası benzerlik özellikleri
  - 1.  $\text{sim}(i,j) \geq 0$
  - 2.  $\text{sim}(i,j) = \text{sim}(j,i)$
-

# İkili Değişkenler Arası Benzerlik

- İkili bir değişkenin 0 veya 1 olarak iki değeri olabilir.
- Bir olasılık tablosu oluşturulur:

		Nesne $j$	
		0	1
Nesne $i$	0	$M_{00}$	$M_{01}$
	1	$M_{10}$	$M_{11}$

$M_{00}$ :  $i$  nesnesinin 0,  $j$  nesnesinin 0 olduğu niteliklerin sayısı

$M_{10}$ :  $i$  nesnesinin 1,  $j$  nesnesinin 0 olduğu niteliklerin sayısı

$M_{01}$ :  $i$  nesnesinin 0,  $j$  nesnesinin 1 olduğu niteliklerin sayısı

$M_{11}$ :  $i$  nesnesinin 1,  $j$  nesnesinin 1 olduğu niteliklerin sayısı

- **Yalın uyum katsayısı** (simple matching coefficient): ikili değişkenin simetrik olduğu durumlarda

$$sim(i, j) = \frac{M_{11} + M_{00}}{M_{00} + M_{01} + M_{10} + M_{11}}$$

- **Jaccard katsayısı** (ikili değişkenin asimetric olduğu durumlar):

$$d(i, j) = \frac{M_{11}}{M_{01} + M_{10} + M_{11}}$$



# Örnek

## □ Jaccard

Name	Fever	Cough	Test-1	Test-2	Test-3	Test-4
Jack	1	0	1	0	0	0
Mary	1	0	1	0	1	0
Jim	1	1	0	0	0	0

$$d(\text{jack}, \text{mary}) = \frac{0 + 1}{2 + 0 + 1} = 0.33$$

$$d(\text{jack}, \text{jim}) = \frac{1 + 1}{1 + 1 + 1} = 0.67$$

$$d(\text{jim}, \text{mary}) = \frac{1 + 2}{1 + 1 + 2} = 0.75$$